

# EXPLOITING DARK INFORMATION RESOURCES TO CREATE NEW VALUE ADDED SERVICES TO STUDY EARTH SCIENCE PHENOMENA

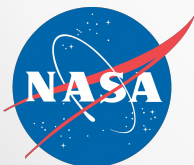
**Rahul Ramachandran** NASA/MSFC, **Manil Maskey**  
UAH, Xiang Li UAH, Kaylin Bugbee UAH

## Project Team:

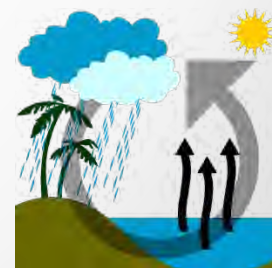
**MSFC/UAH:** Patrick Gatlin, Amanda Weigel, JJ Miller, Ajinkya Kulkarni

**GSFC:** **Steve Kempler**, Suhung Shen, Chung-Lin Shie, Maksym Petrenko

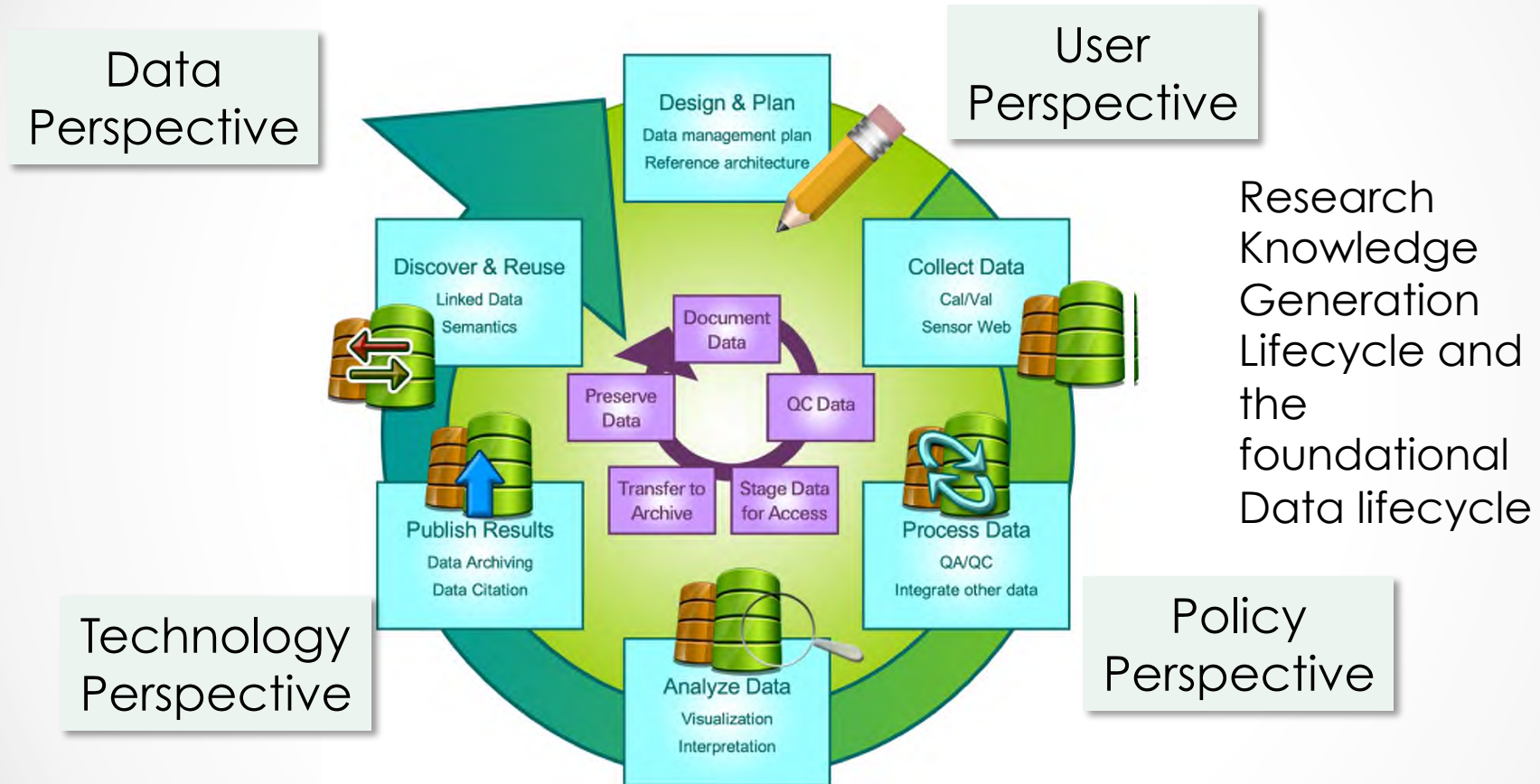
**RPI:** **Peter Fox**, Stefan Zednik, Anirudh Prabhu



**Invited Talk: Earth Observing Data Science**  
**IGARSS July 10-15, 2016 Beijing, China**



# Earth Science Informatics



Goals: to make this process efficient, address existing gaps/hurdles, seamlessly integrate new emerging technology, and enable new research capabilities

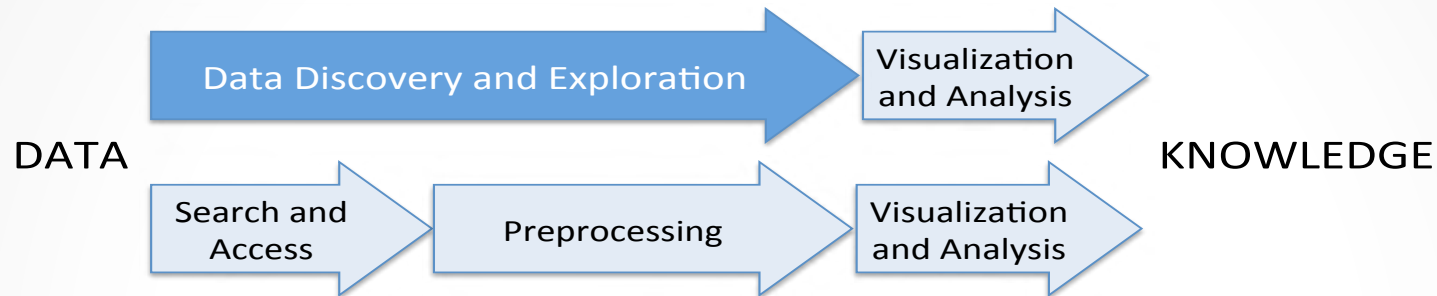
# Outline

1. Project Overview
2. Data Curation Service
3. Rules Engine
4. Application (with Demo)
5. Image Retrieval Service
6. Summary

# Part 1: Project Overview

...

# Motivation



- Data preparation steps are **cumbersome** and **time consuming**
  - Covers discovery, access and preprocessing
- Limitations of current Data/Information Systems
  - **Boolean search** on data based on instrument or geophysical or other **keywords**
  - Underlying **assumption** that users have sufficient knowledge of the **domain vocabulary**
  - **Lack support** for those **unfamiliar** with the domain vocabulary or the **breadth of relevant data** available

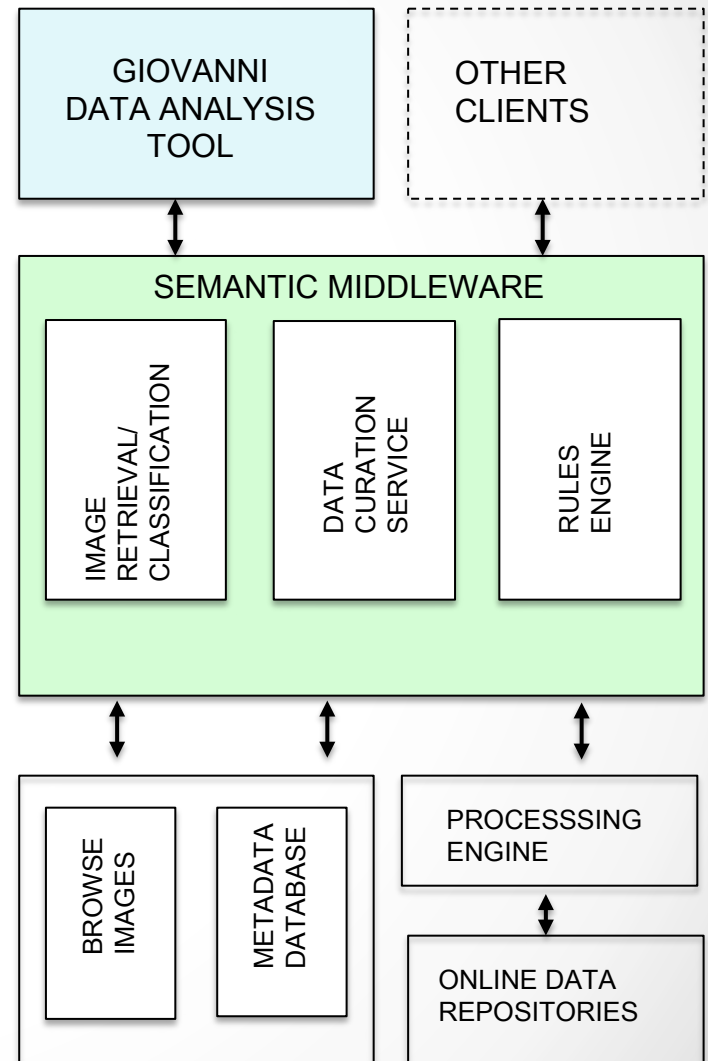
# Earth Science Metadata: Dark Resources

- *Dark resources* - information resources that organizations collect, process, and store for regular business or operational activities but fail to utilize for **other** purposes
  - Challenge is to recognize, identify and effectively utilize these dark data stores
- Metadata catalogs contain dark resources consisting of structured information, free form descriptions of data and browse images.
  - NASA's Common Metadata Repository (CMR) holds >6000 data collections, 270 million records for individual files and 67 million browse images.

Premise: Metadata catalogs can be utilized *beyond their original design intent* to provide *new data discovery and exploration pathways* to support science and education communities.

# Project Goals

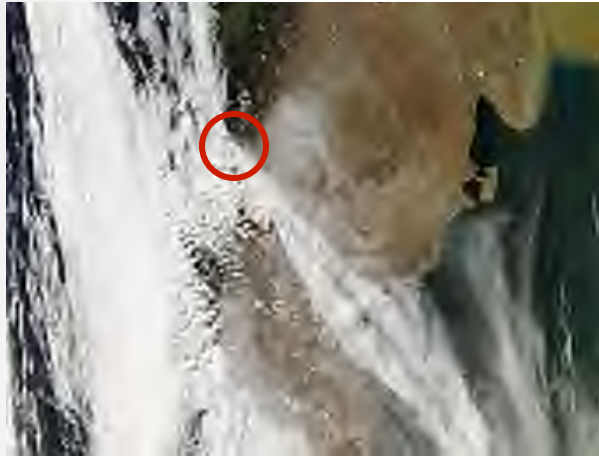
- Design a Semantic Middleware Layer (SML) to exploit these metadata resources
  - provide novel **data discovery and exploration** capabilities that significantly reduce data preparation time.
  - utilize a varied set of semantic web, information retrieval and image mining technologies.
- Design SML as a Service Oriented Architecture (SOA) to allow individual components to be used by existing systems



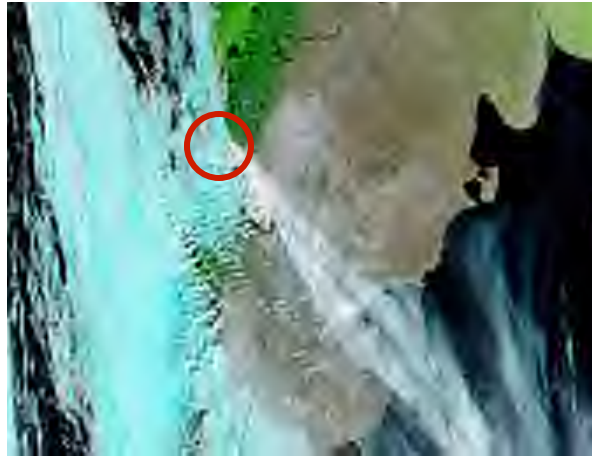


# Use Case:

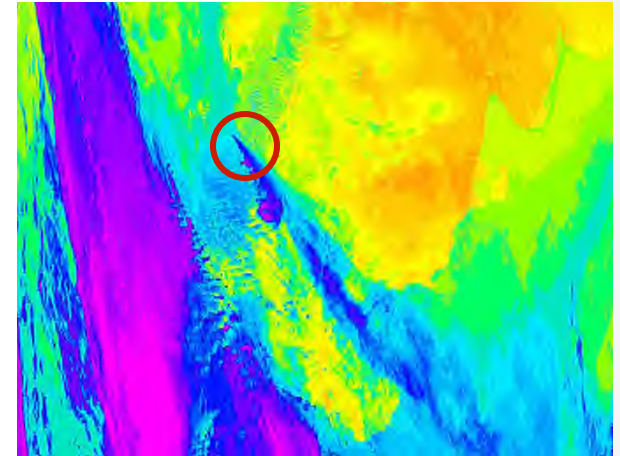
## *Find Interesting Events from Browse Images*



Band 1-4-3 (true color)



Band 7-2-1



LST

Example: MODIS-Aqua 2008-05-03 18:45 UTC

### **Chaitén Volcano Eruption**

**Eruption Time period: May 2 – Nov 2008**

**Location: Andes region, Chile ( -42.832778, -72.645833)**

**Image Retrieval Service can be used to find volcanic ash events in browse imagery**





# Suggest Relevant Data

## Total SO<sub>2</sub> mass:

e.g. **Chaitén** is 10 (kt) =(kilotons ) , (1kt= 1000 metric tons)

[ftp://measures.gsfc.nasa.gov/data/s4pa/SO2/MSVOLSO2L4.1/MSVOLSO2L4\\_v01-00-2014m1002.txt](ftp://measures.gsfc.nasa.gov/data/s4pa/SO2/MSVOLSO2L4.1/MSVOLSO2L4_v01-00-2014m1002.txt)

## Daily SO<sub>2</sub>:

OMI/Aura Sulphur Dioxide (SO<sub>2</sub>) Total Column Daily L2 Global 0.125 deg

[http://disc.sci.gsfc.nasa.gov/datacollection/OMSO2G\\_V003.html](http://disc.sci.gsfc.nasa.gov/datacollection/OMSO2G_V003.html)

## Calibrated Radiances:

MODIS/Aqua Calibrated Radiances 5-Min L1B Swath 1km

<http://dx.doi.org/10.5067/modis/myd021km.006>

## Aerosol Optical Thickness:

MODIS/Aqua Aerosol 5-Min L2 Swath 10km

<http://modis-atmos.gsfc.nasa.gov/MODC>

SeaWiFS Deep Blue Aerosol Optical Depth Data 13.5km

<http://disc.gsfc.nasa.gov/datacollection>

**Data Curation Service  
recommends relevant  
datasets to support event  
analysis**

## IR Brightness Temperature:

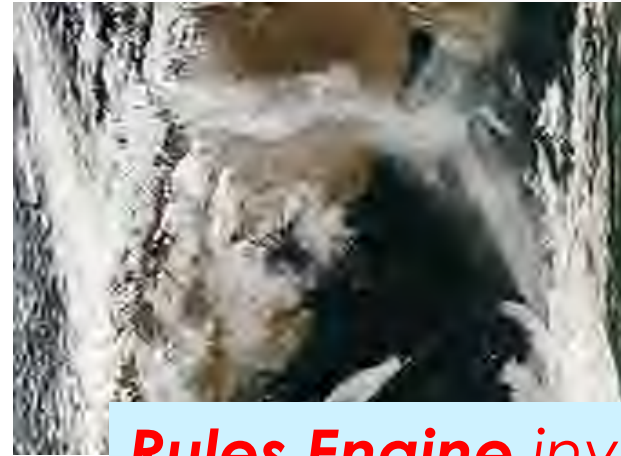
NCEP/CPC 4-km Global (60 deg N - 60 deg S) Merged IR Brightness Temperature Dataset

# Generate Giovanni SO2 Plots

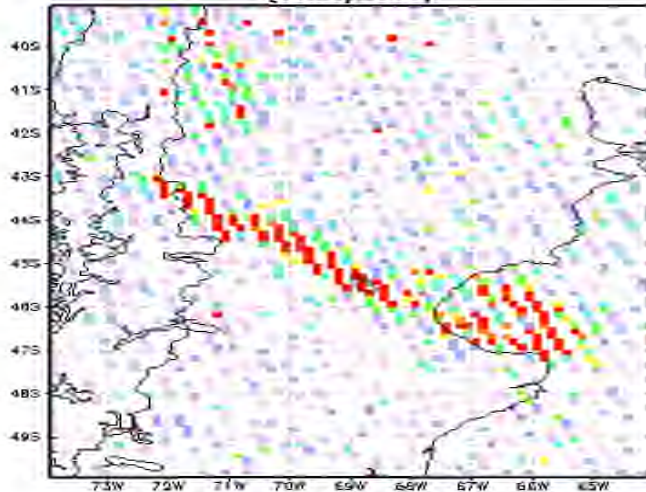
MODIS-Aqua 2008-05-03 18:45 UTC



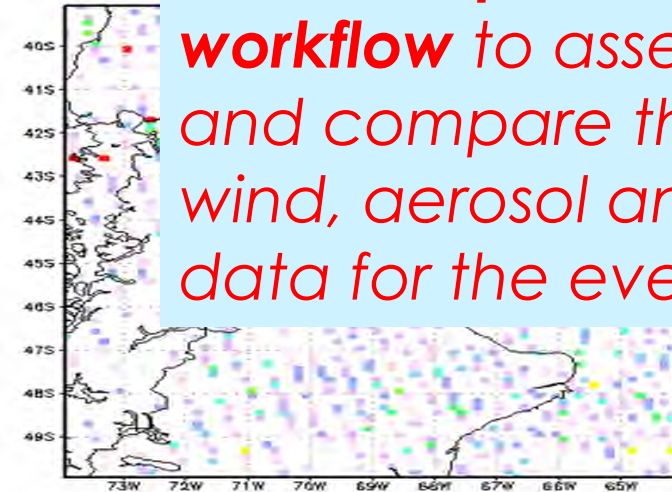
MODIS-Aqua 2008-05-05 18:30 UTC



2G.003 SO2 Column Amount (Planetary Boundary La  
{03May2008})



2G.003 SO2

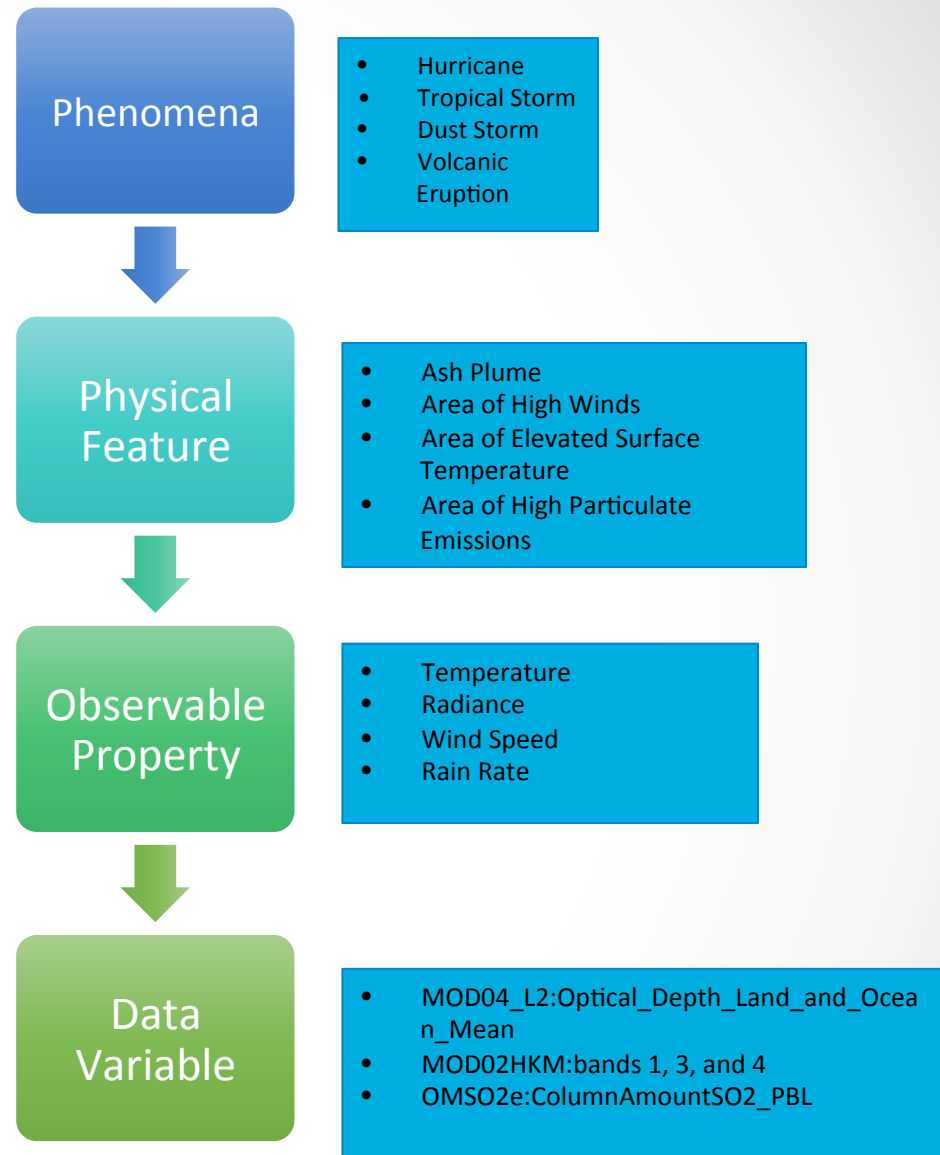


**Rules Engine** invokes a **Giovanni processing workflow** to assemble and compare the wind, aerosol and SO2 data for the event

[http://gdata2.sci.gsfc.nasa.gov/daac-bin/G3/gui.cgi?instance\\_id=omil2g](http://gdata2.sci.gsfc.nasa.gov/daac-bin/G3/gui.cgi?instance_id=omil2g)

# Conceptual Model

- **Phenomena**
  - Event type
- **Physical Feature**
  - Manifestation / Driver of phenomena
  - Has space/time extent
  - Can precede or linger after what is generally thought of as the phenomena event
- **Observable Property**
  - Characteristic/property of physical feature
- **Data Variable**
  - Measurement/estimation of observable feature



# **Part 2: Data Curation**

## **Algorithm for Phenomena**

...

# Data Curation

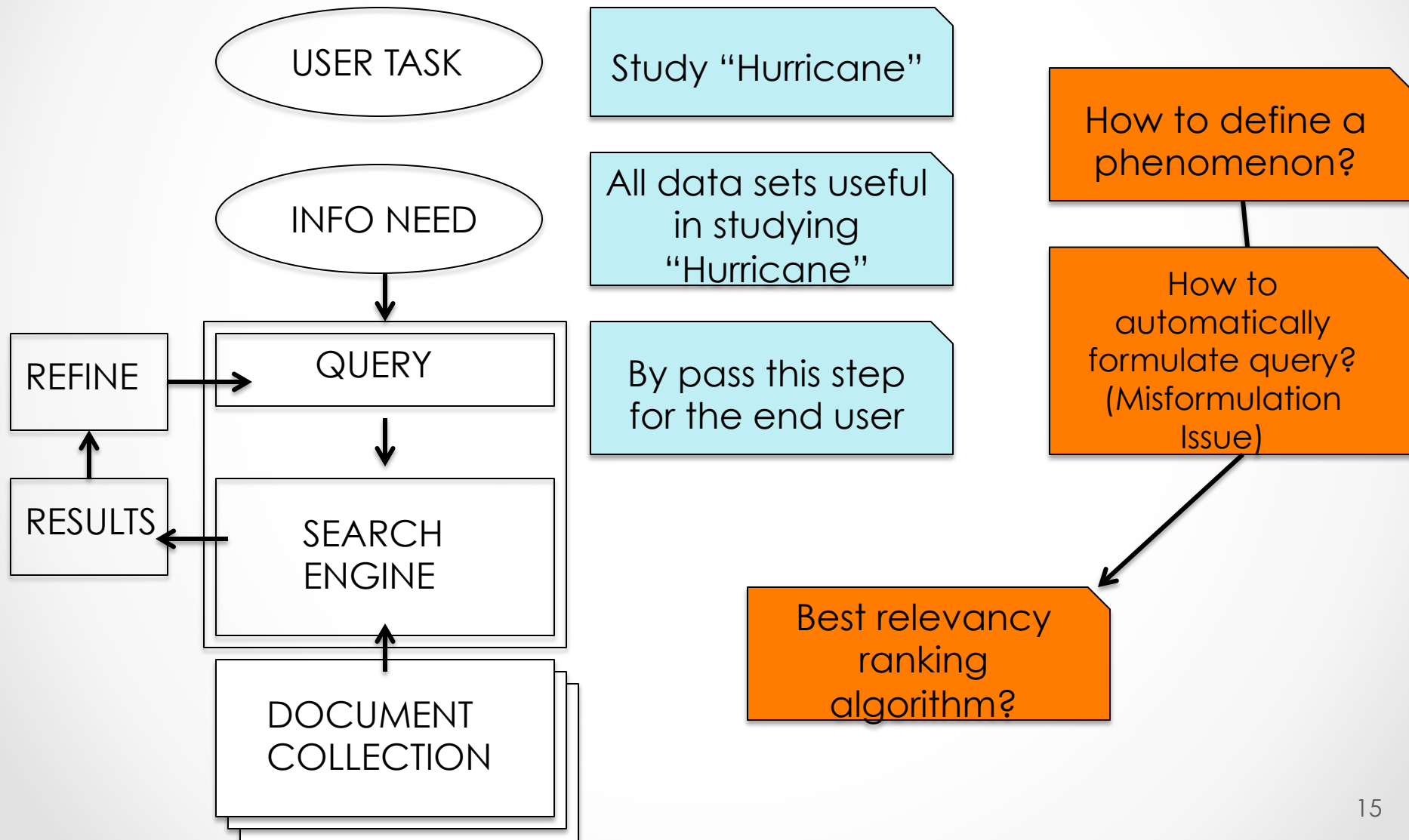
- Curation is traditionally defined as the process of collecting and organizing information around a common subject matter or a topic of interest and typically occurs in museums, art galleries, and libraries.
- Ramachandran et al. [2015] define geocuration as the act of searching, selecting, and synthesizing Earth science data/metadata and information from across disciplines and repositories into a single, cohesive, and useful collection.
  - Manual
  - Automated

# Objectives

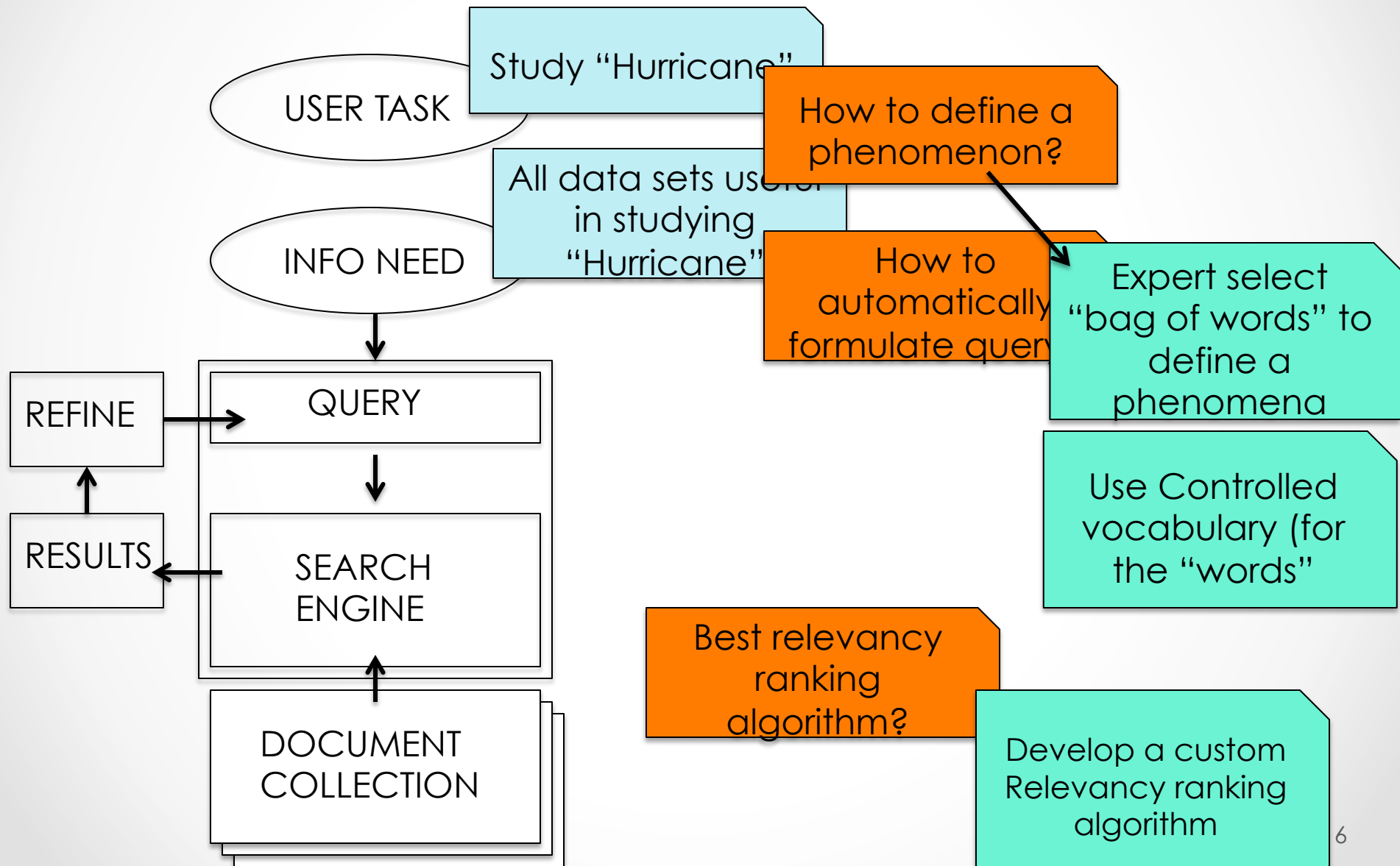
- Design a data curation (relevancy ranking) algorithm for a set of **phenomena**
- Provide the data curation algorithm as a stand alone service
- Envisioned Use:
  - Given a phenomenon type (Ex: Hurricane), DCS returns a list of relevant data sets (variables)
    - $\langle \text{list of data sets (variables)} \rangle = \text{DCS(Phenomenon Type)}$
  - For a specific phenomenon instance (event: Hurricane Katrina), these curated datasets can be filtered based on space/time to get actual granules



# Data Curation is a Specialized Search Problem



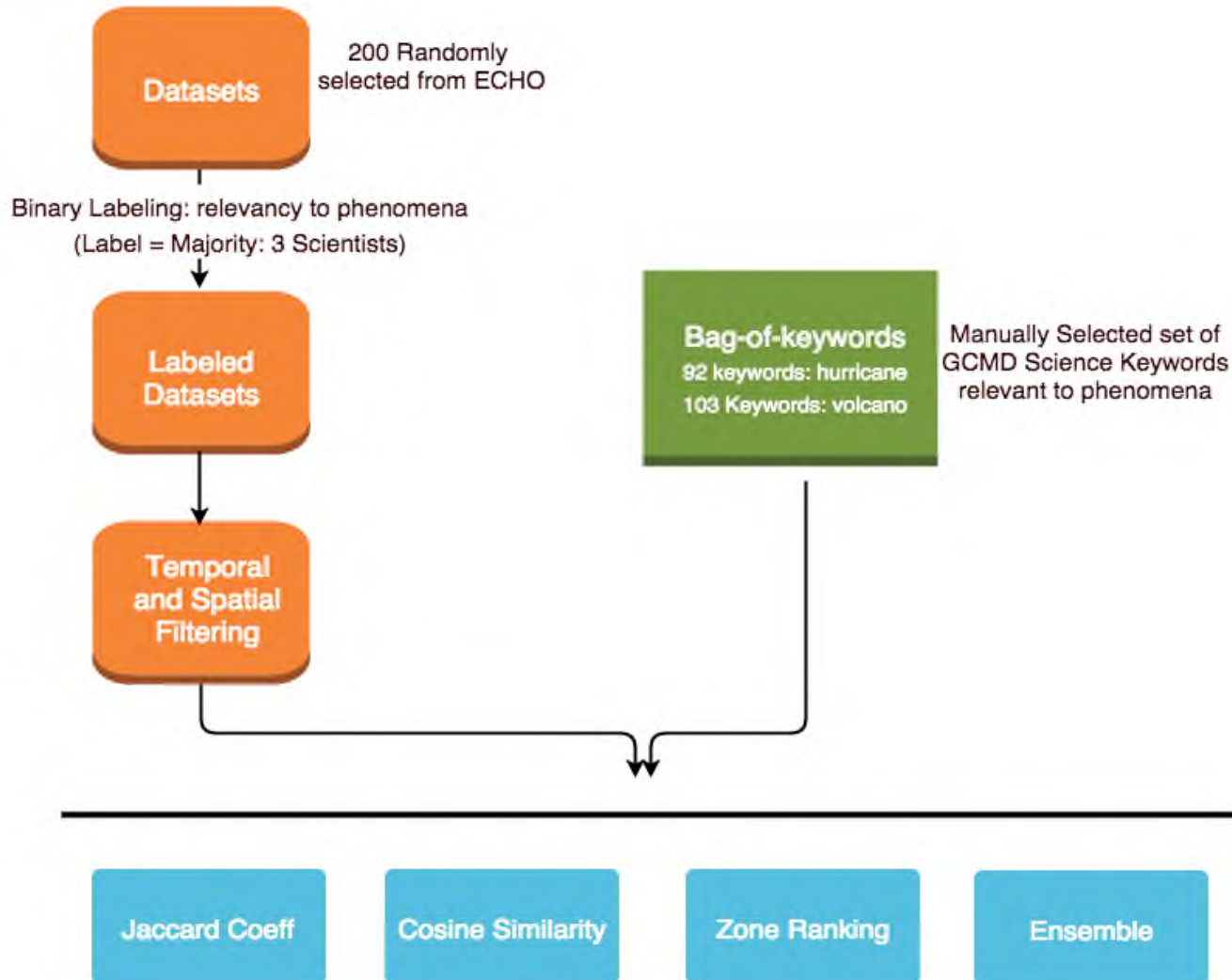
# Our Approach



# Methodology

- Selected three metadata fields
  - Science Keywords
  - Data set long name (title)
  - Data set description
- Developed customized vector space model for each field
- Compared different similarity measures
  - Cosine vs Jaccard
- Used Weighted Zone Ranking (Ensemble)
  - $S_c(e) = w_s \cdot S_c(s) + w_l \cdot S_c(l) + w_d \cdot S_c(d)$

# Experiment Setup



# Comparison of Similarity Measures

	Hurricane		Volcanic Eruption	
	Jaccard Coefficient	Cosine Similarity	Jaccard Coefficient	Cosine Similarity
Top 10 retrieval	10	9	6	7
Top 20 retrieval	17	16	15	15
Top 30 retrieval	23	24	22	21

- Both of the measures performed similarly
- Selected Cosine Similarity measure because it is commonly used in space vector model information retrieval

# Ranking Results (Top 20) using Ensemble Method

	Optimal Weight		Equal Weight		Random	
	Precision	Recall	Precision	Recall	Precision	Recall
Hurricane	90.0%	47.4%	85.0%	44.7%	54.3%	28.6%
Volcanic Eruption	85.0%	68.0%	80.0%	64.0%	62.5%	50.0%
Fire	75.0%	30.0%	75.0%	30.0%	64.1%	25.6%
Flood	65.0%	48.1%	55.0%	40.7%	35.5%	26.3%

- Different numbers of “relevant” data sets, collection size (recall) exist within each truth set for each phenomenon
- Better to compare the curation results against the random selection rather than compare the performance against each other
- On average, precision improves about 25% when using our method and recall improves about 16%



# Optimal ensemble weights for each phenomenon

Phenomenon	Optimal Weight Set ( $W_{\text{sciencekeyword}}$ , $W_{\text{longname}}$ , $W_{\text{description}}$ )
Hurricane	(0.6, 0.1, 0.3)
Volcanic Eruption	(0.2, 0.6, 0.2)
Fire	(0.6, 0.2, 0.2)
Flood	(0.5, 0.4, 0.1)

- Weight for science keyword is largest while the weight for description is smallest
  - Science keywords metadata fields use a controlled vocabulary and should be accurate and consistent
  - Description field is free-text and has the most variability in quality

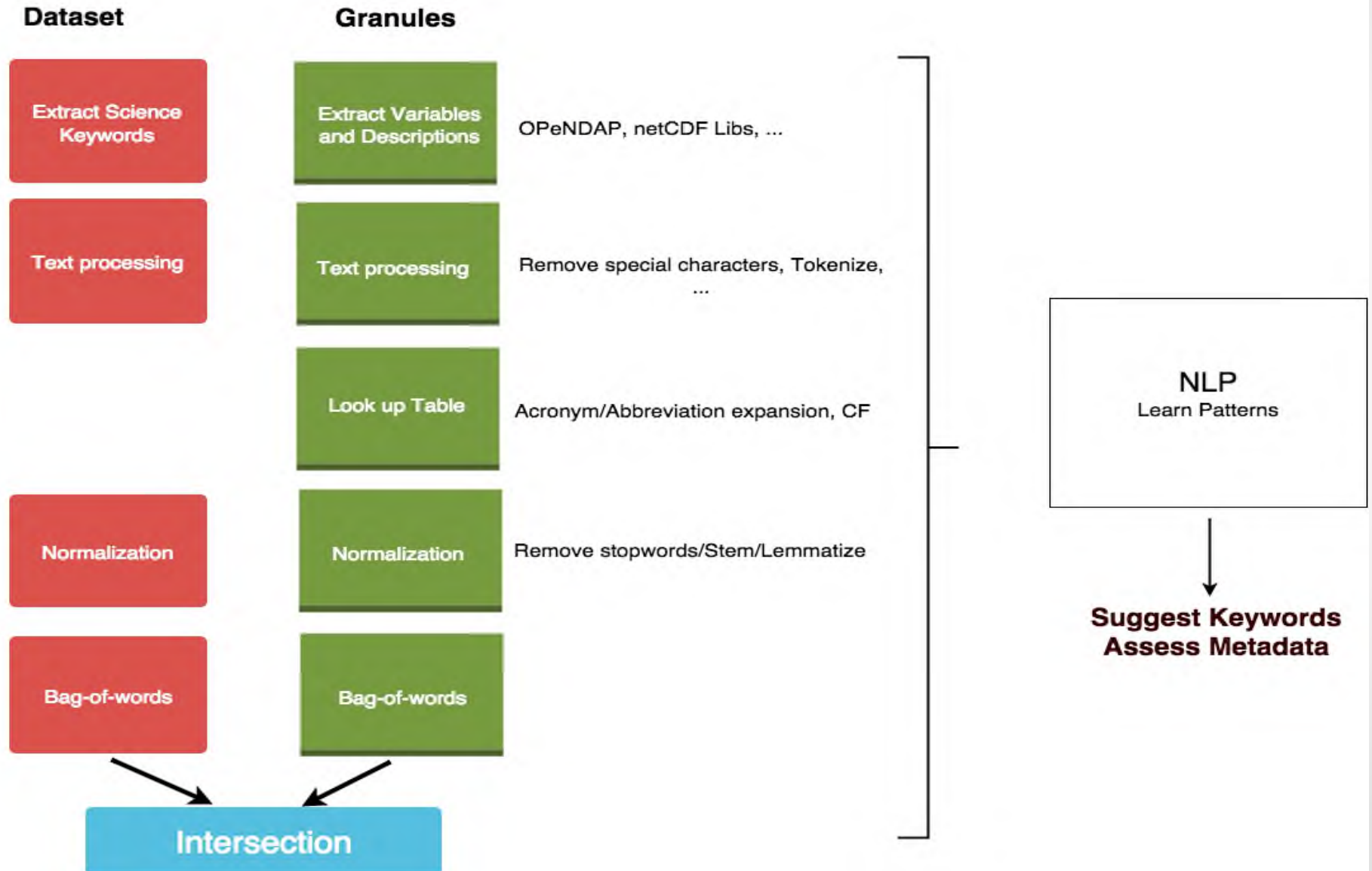
# Methodology Limitations

- *Modeling the search intent is difficult*
  - one may be interested in only a specific aspect of a phenomenon whereas another user may only be interested in some other characteristic of a phenomenon
- *Quality of metadata records is variable*
  - Key assumption is that the metadata records stored in the CMR catalog are consistent, correct, and complete
  - Launched a project to fix this
- *Granularity of the Controlled Vocabulary*
  - Rich detailed controlled vocabulary provides a better level of annotation granularity to represent different phenomena and help disambiguate data sets
- *Truth set labels may be biased*
  - domain experts on our team have stronger expertise in certain areas such as hurricanes and weaker expertise in others

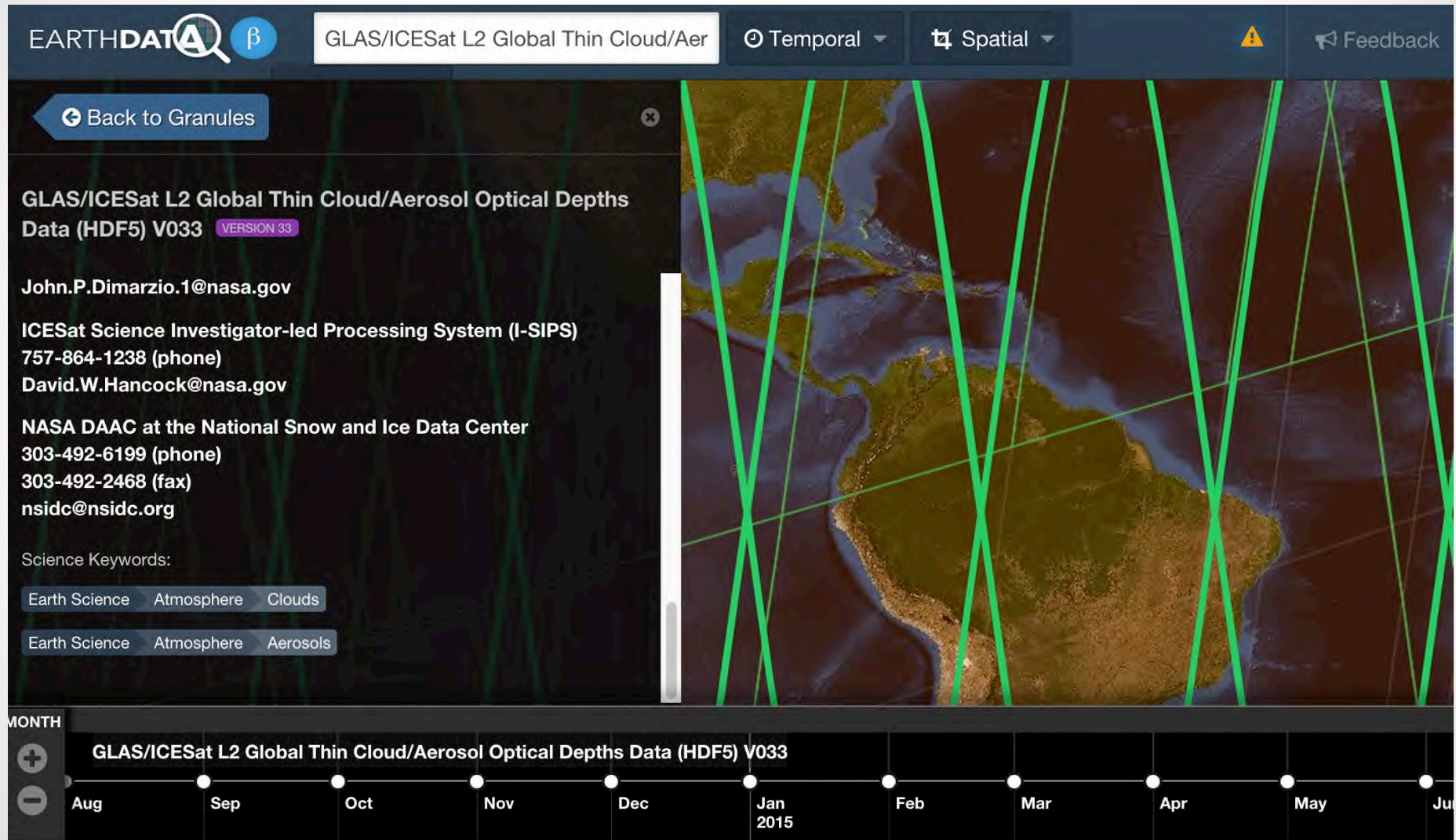
# Next: Find relevant data fields

- Need actual data variables
  - Example: Giovanni uses these fields for visualization
- What we know
  - Data set (Collection) level science keywords (GCMD) – Experts
  - Granule data fields and metadata – Auto extract\*
- How do we map?
  - Start with GCMD to CF Standard name
  - Most don't follow CF Standard names

# Approach

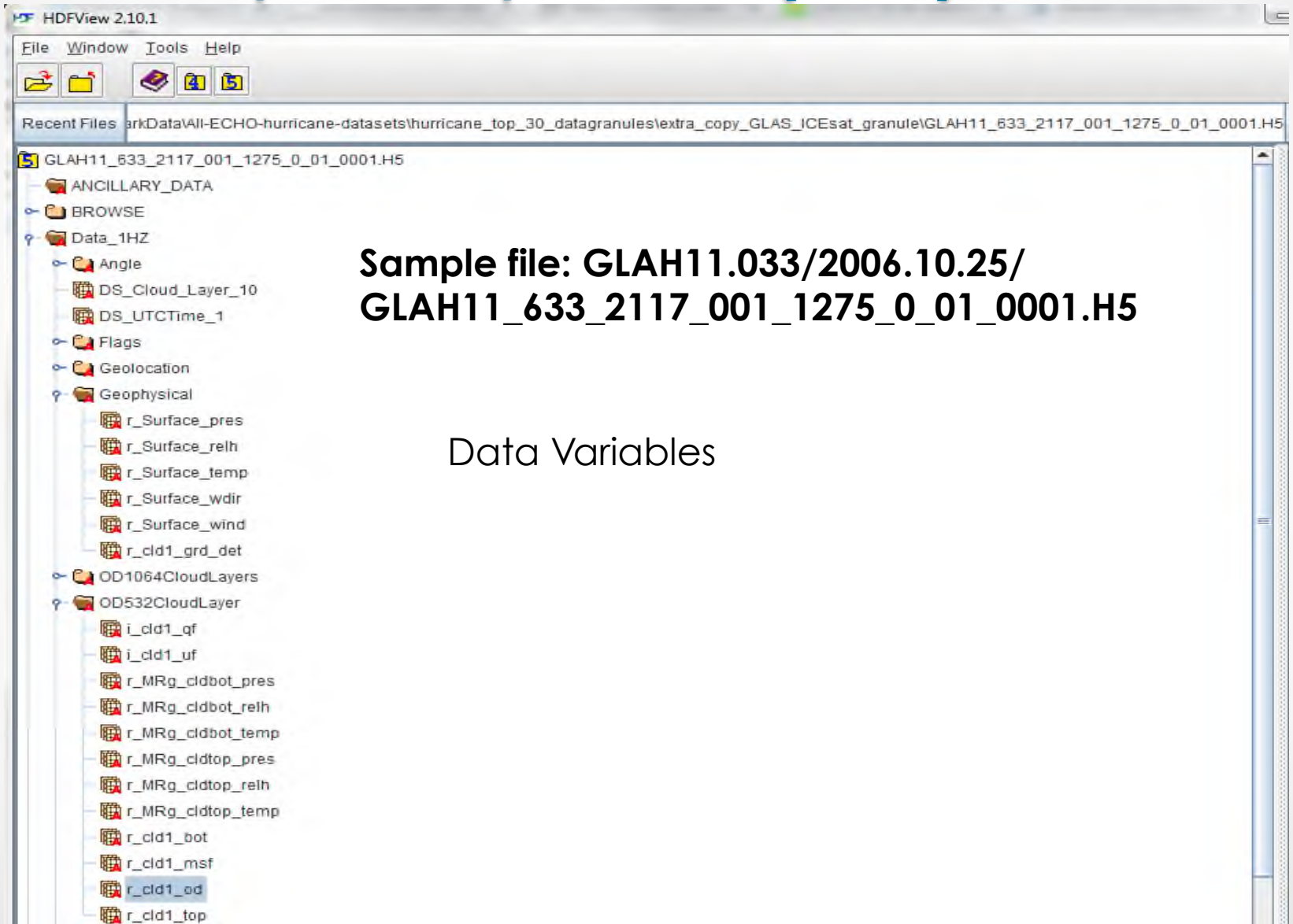


# Example: GLAS/ICESat L2 Global Thin Cloud/Aerosol Optical Depths Data (HDF5) V033 – Dataset Metadata





# Example: GLAS/ICESat L2 Global Thin Cloud/Aerosol Optical Depths Data (HDF5) V033





# Example: GLASICESat L2 Global Thin Cloud Aerosol Optical Depths Data (HDF5) V033

## Science keyword to variable mapping

- r\_Surface\_relh | Surface Relative Humidity
  - No match
- r\_Surface\_temp | Surface Temperature
  - No match
- r\_Surface\_wind | Surface Wind Speed
  - No match
- r\_cld1\_od | Cloud Optical Depth at 532 nm
  - Score=3 keyword: ATMOSPHERE->CLOUDS->CLOUD OPTICAL DEPTH/THICKNESS
  - Score=2 keyword: ATMOSPHERE->AEROSOLS->AEROSOL OPTICAL DEPTH/THICKNESS

## Variable to keyword mapping

- ATMOSPHERE->CLOUDS->CLOUD OPTICAL DEPTH/THICKNESS
  - Score=3 name: r\_cld\_ir\_OD | Cloud Optical Depth at 1064 nm
  - score=3 name: i\_cld1\_qf | Cloud optical depth flag for 532 nm
  - Score=3 name: i\_cld1\_uf | Cloud optical depth flag for 532 nm
  - Score=3 name: r\_cld1\_od | Cloud Optical Depth at 532 nm
  - more with low scores

- **Serendipitous Discovery** - Data Curation Parameter Mapping Algorithm can be used to assess
  - Metadata quality for both dataset and granules
  - Find incorrect/incomplete keyword annotations
  - Automatically suggest science keywords

# Parameter Mapping Tool

54.172.157.10:5000

News personal Mendeley GKeep NASA Demo HS3 GHRC DarkData nspires RResp Unisys Weather - Ge

## Data Parameter Mapping Tool

### Datasets

AIRS/Aqua Level 2 Support retrieval (AIRS+AMSU) V005
GHRST Level 2P USA NASA MODIS Aqua SST:1
MODIS/Terra Temperature and Water Vapor Profiles 5-Min L2 Swath 5km V005
LIS/OTD 2.5 DEGREE LOW RESOLUTION DIURNAL CLIMATOLOGY (LRDC) V2.3.2013
MODIS/Terra Aerosol 5-Min L2 Swath 10km V005 NRT

**Datasets**

MODIS/Terra Aerosol 5-Min L2 Swath 10km V005 NRT

### Science Keyword Map

**EDIT**

ATMOSPHERE → AEROSOLS → PARTICULATE_MATTER 1
Deep_Blue_Aerosol_Optical_Depth_Land_STD : 1
Deep_Blue_Aerosol_Optical_Depth_550_Land : 1
Aerosol_Type_Land : 1
Aerosol_Cldmask_Byproducts_Ocean : 1
Deep_Blue_Aerosol_Optical_Depth_Land : 1
Aerosol_Cldmask_Byproducts_Land : 1
Deep_Blue_Aerosol_Optical_Depth_550_Land_STD : 1
Optical_Depth_Small_Average_Ocean : 0

### Parameter Map

**EDIT**

Optical_Depth_Small_Average_Ocean 1
ATMOSPHERE → AEROSOLS → AEROSOLS_OPTICAL_DEPTH_THICKNESS : 2
ATMOSPHERE → ATMOSPHERIC_RADIATION → OPTICAL_DEPTH_THICKNESS : 2
ATMOSPHERE → AEROSOLS → PARTICULATE_MATTER : 0
Asymmetry_Factor_Best_Ocean 0
Deep_Blue_Angstrom_Exponent_Land 0
Cloud_Fraction_Ocean 1
ATMOSPHERE → AEROSOLS → CLOUD_CONDENSATION_NUCLEI : 1

**Science Keyword**

**Parameter**

**Mapped Science Keywords**

**Mapped Parameters**

MODIS/Terra Aerosol 5-Min L2 Swath 10km V005 NRT

ATMOSPHERE → AEROSOLS → AEROSOL_PARTICLE_PROPERTIES : 2	Remove
ATMOSPHERE → AEROSOLS → CLOUD_CONDENSATION_NUCLEI : 2	Remove
ATMOSPHERE → AEROSOLS → AEROSOL_EXTINCTION : 2	Remove
ATMOSPHERE → AEROSOLS → AEROSOLS_OPTICAL_DEPTH_THICKNESS : 2	Remove
ATMOSPHERE → AEROSOLS → AEROSOL_RADIANCE : 2	Remove
ATMOSPHERE → AEROSOLS → CARBONACEOUS_AEROSOLS : 2	Remove
ATMOSPHERE → AEROSOLS → DUST/ASH/SMOKE : 2	Remove
ATMOSPHERE → AEROSOLS → NITRATE_PARTICLES : 2	Remove
ATMOSPHERE → AEROSOLS → ORGANIC_PARTICLES : 2	Remove
ATMOSPHERE → AEROSOLS → PARTICULATE_MATTER : 2	Remove
ATMOSPHERE → AEROSOLS → SULFATE_PARTICLES : 2	Remove
ATMOSPHERE → ATMOSPHERIC_RADIATION → RADIATIVE_FLUX : 2	Remove
ATMOSPHERE → ATMOSPHERIC_RADIATION → REFLECTANCE : 2	Remove
✓ ATMOSPHERE → ATMOSPHERIC_RADIATION → OPTICAL_DEPTH_THICKNESS : 2	Remove

ATMOSPHERE → AEROSOLS → PARTICULATE\_MATTER : 0

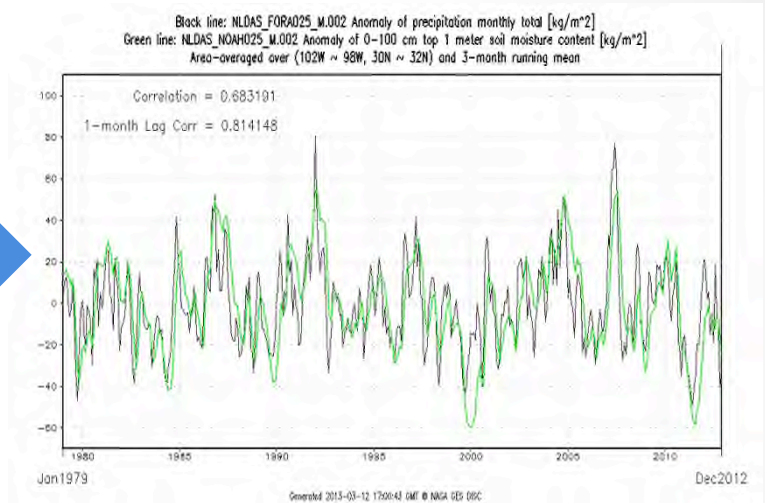
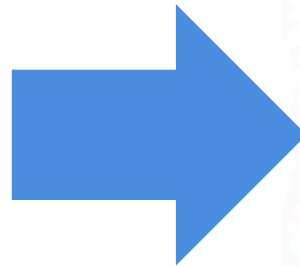
Remove

**Edit/Save Mapping**

**Mapping Scores Generated by Algorithm**

# **Part 3: Rules Engine**

# What settings should I use to visualize this event?



Data  
Variable  
?

Dataset  
?  
Visualization  
Type?

Goal: Automate data preprocessing and exploratory analysis and visualization tasks

# Strategy

- Service to generate and rank candidate workflow configurations
- Use rules to make **assertions** about **compatibility based on multiple factors**
  - does this data variable make sense for this feature?
  - does this visualization type make sense for this feature?
  - does the temporal / spatial resolution of this dataset make sense for this feature?
- Each compatibility assertion type is assigned weights.
  - ex: Strong = 5, Some = 3, Slight = 1, Indifferent = 0, Negative = -1.
- Based on the aggregated compatibility assertions, we calculate the score for each visualization candidate.

# Ruleset Development

Survey asked users to rate characteristics of phenomena features

## Feature characteristics for analysis \*

What characteristics are of interest when analyzing the feature?

	negative value	indifferent	slight value	some value	strong value
east-west movement	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
north-south movement	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
temporal evolution	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
spatial extent of event	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
year-to-year variability	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
may impact seasonal variation	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
variation with atmospheric height	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
global phenomena	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
detection of events	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>



Survey results used to formulate rules

[rule1:

(?feature rdf:type  
dd:AshPlume)

->

(?feature  
dd:strongCompatibilityFor  
dd:temporal\_evolution),

(?feature  
dd:indifferentCompatibilityFor  
dd:east-west-movement),

...

]



# Phenomena Feature Characteristic Mappings

Phenomena	East-West Movement	North-South Movement	Temporal Evolution	Spatial Extent of Event	Year-to-Year Variability	May Impact Seasonal Variation	Variation with Atmospheric Height	Global Phenomena	Detection of Events
<b>Volcano - Ash Plume</b>	Indifferent	Indifferent	Strong	Slight	Strong	Strong	Strong	Strong	Strong
<b>Flood</b>	Some	Some	Strong	Some	Some	Strong	Some	Slight	Some
<b>Dust Storm</b>	Strong	Strong	Strong	Strong	Indifferent	Indifferent	Strong	Indifferent	Some

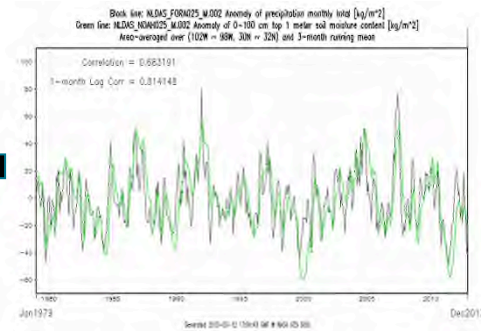
# Service to Characteristic Mappings

Service	Visualization	East-West Movement	North-South Movement	Temporal Evolution	Spatial Extent of Event	Year-to-Year Variability	Seasonal Variation	Variation with Atmospheric Height	Global Phenomena	Detection of Events
Time-averaged Map	Color-Slice Map				✓					
Area-averaged Time Series	Time Series			✓						✓
User-defined Climatology	Color-Slice Map						✓			
Vertical Profile	Line Plot							✓		
Seasonal Time Series	Time Series					✓				
Zonal Means	Line Plot								✓	
Hovmoller (Longitude)	Color-Slice Grid	✓								
Hovmoller (Latitude)	Color-Slice Grid		✓							

# Compute Compatibility



+



=



Phenomena:  
Volcano - Ash  
Plume

Service - Area  
Averaged Time  
Series

**STRONG  
COMPATIBILITY  
x2**

Temporal Evolution	Detection of Events
<b>Strong</b>	<b>Strong</b>

Area Averaged Time Series : bestFor →	Temporal evolution; Detection of events
---------------------------------------	---

Images from , [http://disc.sci.gsfc.nasa.gov/data/releases/images/nldas\\_monthly\\_climatology\\_figure\\_9.gif](http://disc.sci.gsfc.nasa.gov/data/releases/images/nldas_monthly_climatology_figure_9.gif), <http://www.clipartbest.com/cliparts/biy/bAX/biybAXGIL.png>

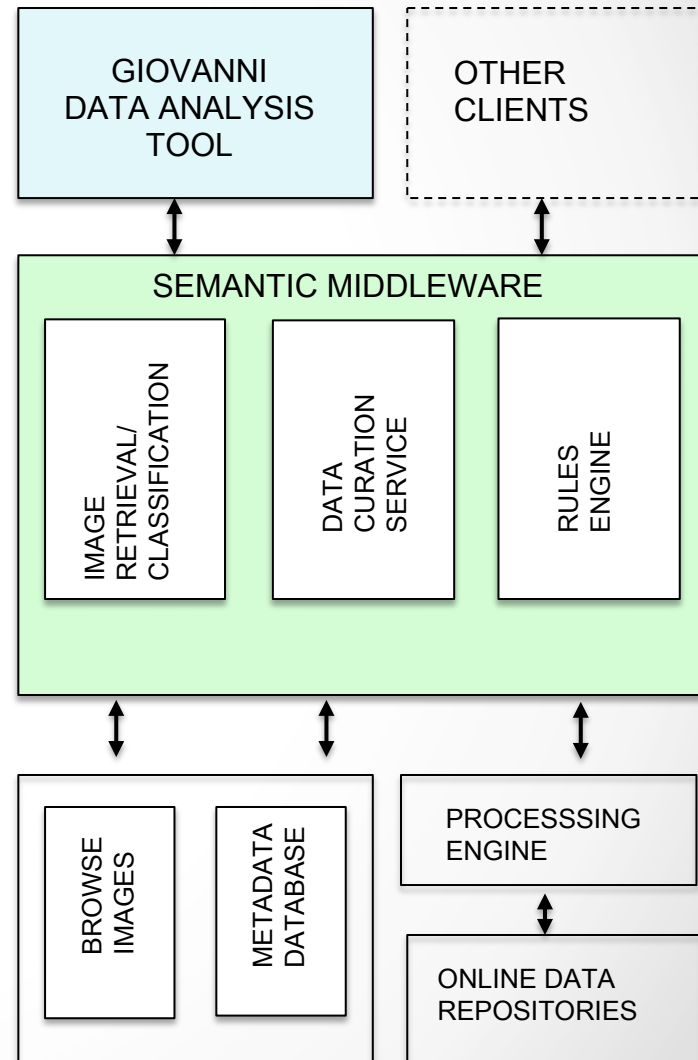
volcanic ash image - By Boaworm (Own work) [CC BY 3.0 (<http://creativecommons.org/licenses/by/3.0>)], via Wikimedia Commons

# **Part 4: Application (Demo)**

# Integrating Services in Giovanni

- **Tool:** Giovanni is a popular on-line environment that lets users discover, plot, and download a number of geophysical parameters (data variables)
- **Goal:** Leverage Dark Data services and technologies to assist Giovanni users in discovering and exploring data

*'Success will be realized when Giovanni requests can be automatically invoked with the appropriate spatial and temporal extents, variables and workflow / visualization type for a particular event'*



# Giovanni – Standard Edition

The screenshot displays the Giovanni Standard Edition web interface. At the top, the header includes 'EARTHDATA', 'Data Discovery', 'Data Center', 'Community', and 'Science Disciplines'. The main title 'GIOVANNI' is followed by the tagline 'The Bridge Between Data and Science' and version 'v 4.17.2'. Below this, there are links for 'Release Notes', 'Browser Compatibility', and 'Known Issues'.

The 'Select Plot' section features several dropdown menus: 'Maps: Time Averaged Map' (selected), 'Comparisons: Select...', 'Time Series: Select...', 'Vertical: Select...', and 'Miscellaneous: Select...'. The 'Select Date Range (UTC)' section includes input fields for 'YYYYMMDD' and 'MM/DD/YYYY', with a 'Valid Range: 1979-01-01 to 2016-02-04' note. The 'Select Region (Bounding Box or Point)' section has a 'Format: Lat, Lon, Bbox, Point' dropdown and a text input field showing '-180,-90,180,90'.

The 'Select Variables' section on the left lists 'Disciplines' and 'Measurements' with checkboxes. The 'Number of matching Variables: 721 of 975' and 'Total Variable(s) Included' are displayed. A 'Keyword:' search bar is present. The main table lists variables with columns for 'Variable Name', 'Source', 'Frequency', 'Resolution', 'Start Date', and 'End Date'. A 'Vertical Choices' modal window is open, showing options for 'Cross Map, Latitude-Pressure', 'Cross Map, Longitude-Pressure', 'Cross Map, Time-Pressure', and 'Vertical Profile'.

At the bottom, there are 'Help', 'Reset', 'Feedback', and 'Plot Data' buttons.

User needs to decide:

- Variable(s)
- Time
- Space
- Plot type

<http://giovanni.sci.gsfc.nasa.gov/giovanni/>



# Giovanni – Dark Data Edition

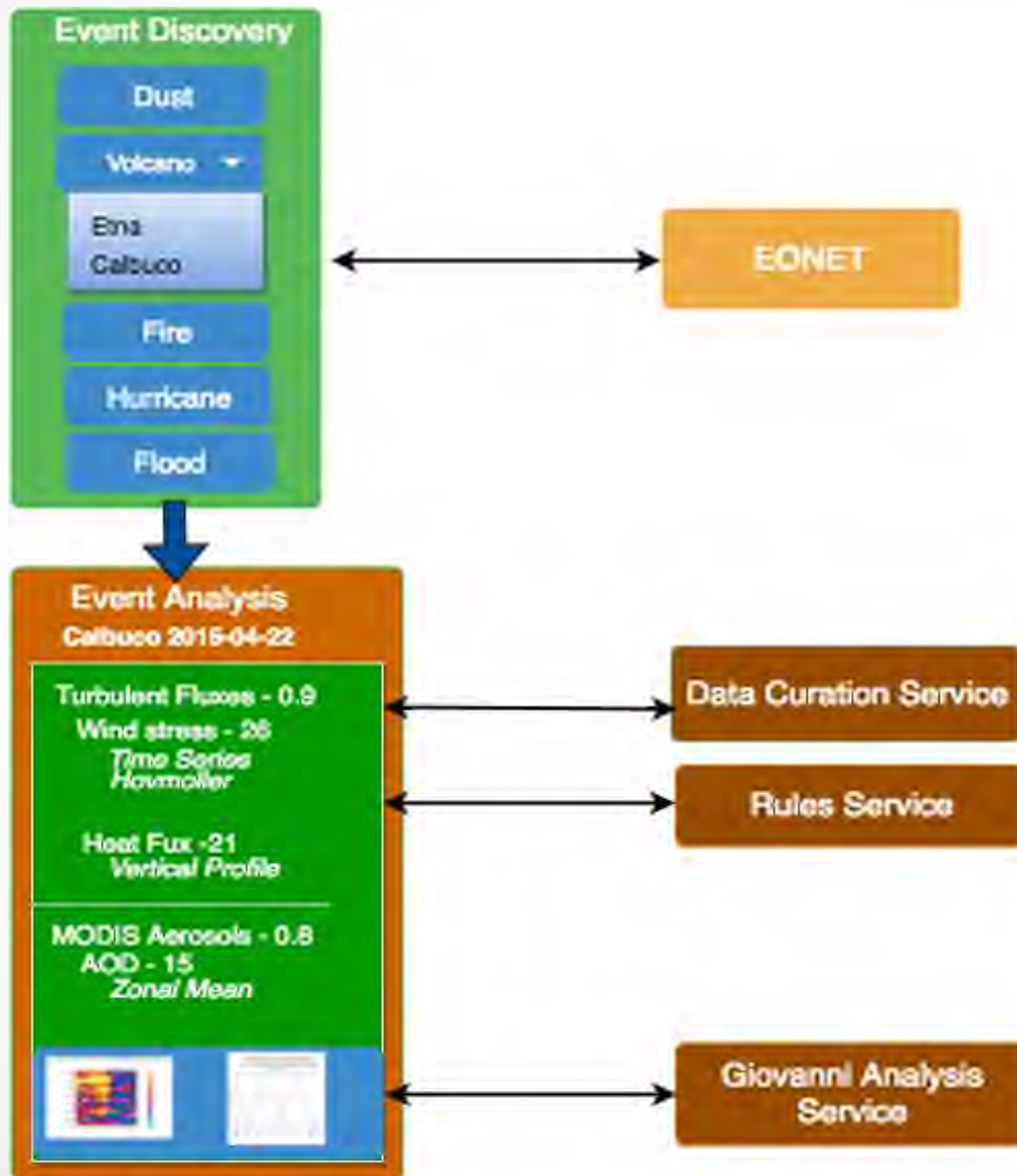
Rules Service:  
highlights  
suitable plots  
based on  
selected event  
& variables

Curation  
Service: event  
type filters  
relevant  
variables

Selected event & its time Event Client

The screenshot shows the GIOVANNI web interface. The top navigation bar includes links for Data Discovery, DAACs, Community, and Scientific Disciplines. The main header reads "GIOVANNI The Bridge Between Data and Science v 4.18" with links for Release Notes, Browser Compatibility, and Known Issues. The "Select Plot" section has a dropdown menu set to "Time Series" and a "Miscellaneous" dropdown. The "Select Date Range (UTC)" section shows a date range from 2015-07-01 to 2015-11-24. The "Select Region (Bounding Box or Shapefile or Event)" section has a dropdown set to "Volcanoes: Manam Volcano". The "Select Variables" section has a sidebar with "Events (all)" selected, showing "Hurricane (14)" and "Volcano (14)". The "Events (by products)" section shows "Hurricane (2)" and "Volcano (2)". The "Events (by variables)" section shows "Hurricane (14)" and "Volcano (14)". The "Disciplines" section lists various categories like Measurements, Platform / Instrument, Spatial Resolutions, Temporal Resolutions, Wavelengths, and Portal. The "Variable" list includes "Aerosol Optical Depth (MY)", "Aerosol Optical Depth (MY)", "Precipitable Water (MY)", "Total Column Water (MY)", "Cirrus Reflected Mean (MY)", "Ice Cloud Optical Mean (MY)", "Liquid Water Mean of Day (MY)", "Cloud Top Pressure (MY)", "Cloud Top Temperature (MY)", "Cloud Top Humidity (MY)", "Cloud Top Albedo (MY)", "Cloud Top Emissivity (MY)", "Cloud Top Reflectance (MY)", "Cloud Top Transmittance (MY)", "Cloud Top Absorptance (MY)", "Cloud Top Emissivity (MY)", "Cloud Top Reflectance (MY)", "Cloud Top Transmittance (MY)", "Cloud Top Absorptance (MY)". The "Event" list includes "Calbuco Volcano, Chile", "Cotopaxi Volcano, Ecuador", "Manam Volcano", "Masaya Volcano, Nicaragua", "Momotombo Volcano, Nicaragua", "Mount Etna Volcano, Italy", "Raging Volcano, Indonesia, July-Aug 2014". The "Event Client" section shows the "Manam Volcano" event selected. The bottom of the interface has buttons for Help, Reset, Feedback, Plot Data, and Go to Results.

# Giovanni - Dark Data Edition



Event  
Analysis  
Workflow

# DEMO



# Part 5: Image Retrieval

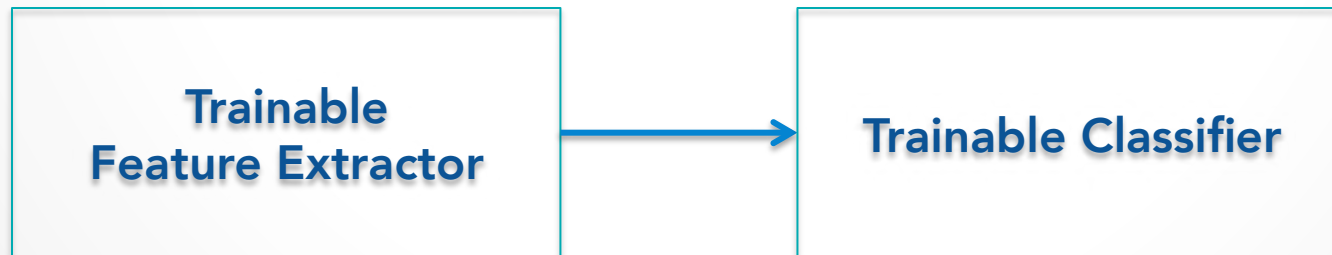
...

# Image Retrieval

- Goal: given an image of Earth science phenomenon retrieve similar images
- Challenge: “semantic gap”
  - low-level image pixels and high-level semantic concepts perceived by humans

# “Deep” Architecture

- Features are key to recognition
- What about learning the features?
- Deep Learning
  - Hierarchical Learning
  - Mimics the human brain that is organized in a deep architecture
  - Processes information through multiple stages of transformation and representation



***Convolutional Neural Network (CNN) - Applicable to Images***



# Transfer Learning

- CNN requires large number of parameters
- Learning parameters from *a few thousand training samples* is unrealistic
- Transfer learning
  - Use internal representation learned from one classification task to another
    - AlexNet architecture - Krizhevsky et. al.
    - Weights learned from ImageNet 1.3 million high-resolution images
    - State-of-the-art classification accuracy

# Experiment: CNN Configuration

Text

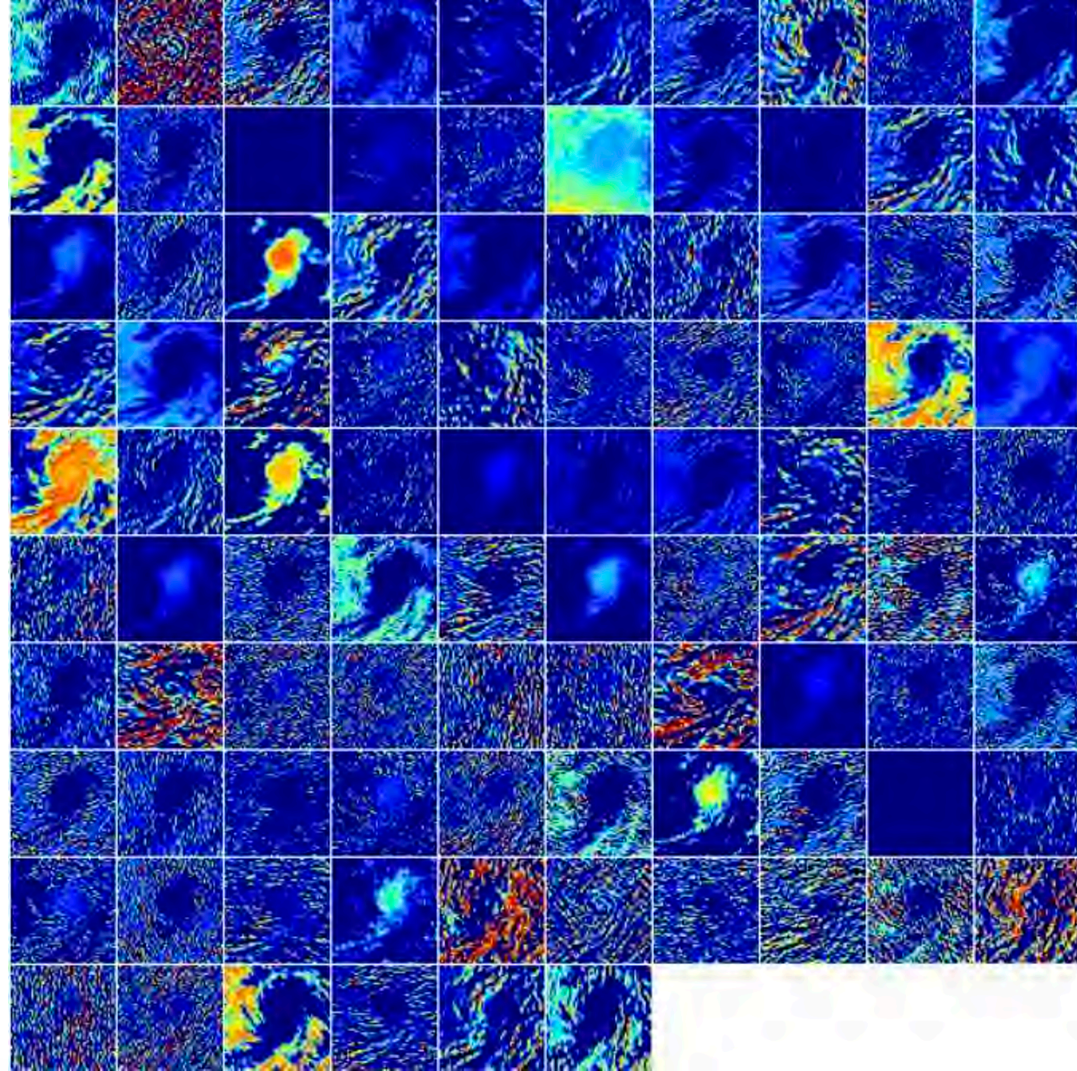


- AlexNet architecture
  - Initialized weights with ImageNet trained model
  - Adaptive learning rate
  - GPU implementation

# Experiment CNN – Visualization



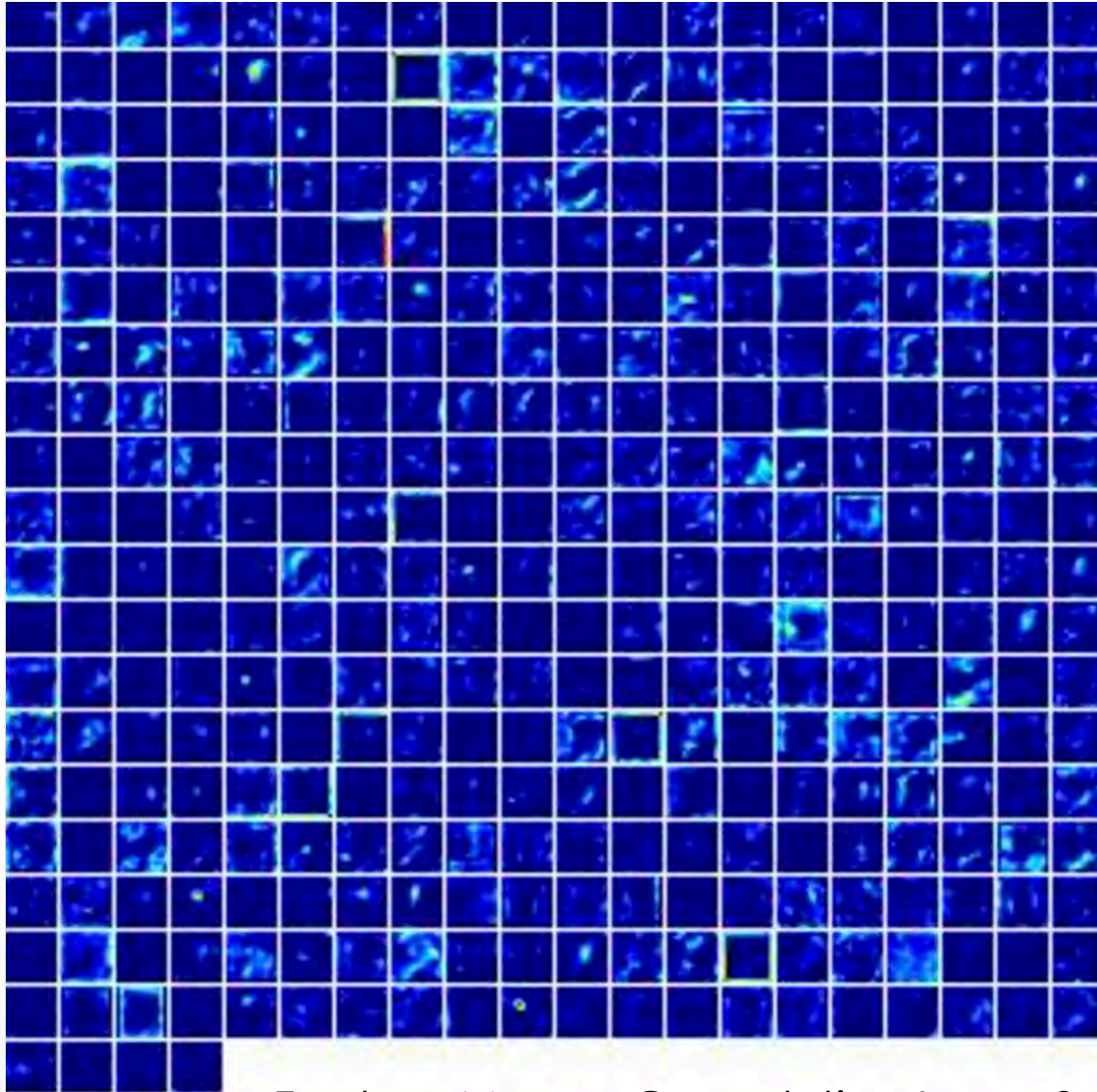
Input Image



Feature Maps – Convolution Layer 1



# Experiment CNN – Visualization



Feature Maps – Convolution Layer 3

# Results: Confusion Matrix

MODIS Rapid Response Test Images (Images are New to Trained CNN)

True/Pred	Dust	Hurricane	Smoke	Other
Dust	<b>287</b>	8	32	33
Hurricane	0	<b>379</b>	1	10
Smoke	12	12	<b>443</b>	9
Other	33	9	23	<b>211</b>

Overall Accuracy = **87.88%**

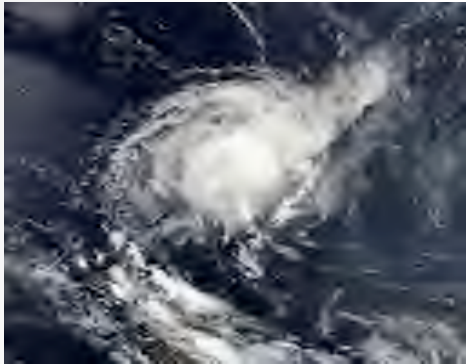
## Producer's Accuracy

Dust 86.45%  
Hurricane 92.89%  
Smoke 88.78%  
Other 80.23%

## User's Accuracy

Dust 79.72%  
Hurricane 97.18%  
Smoke 93.07%  
Other 76.45%

# Results (MODIS Rapid Response)



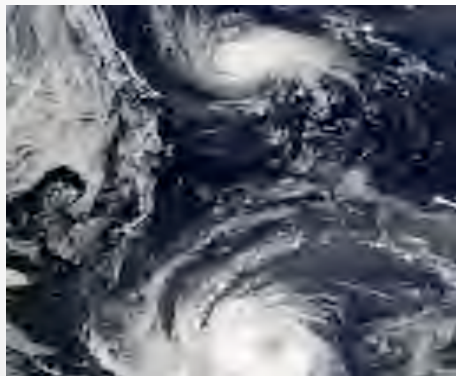
Hurricane – True Positive



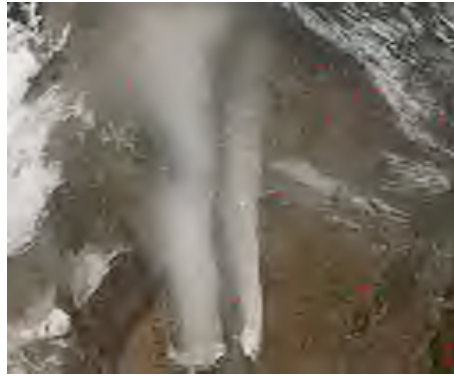
Dust – True Positive



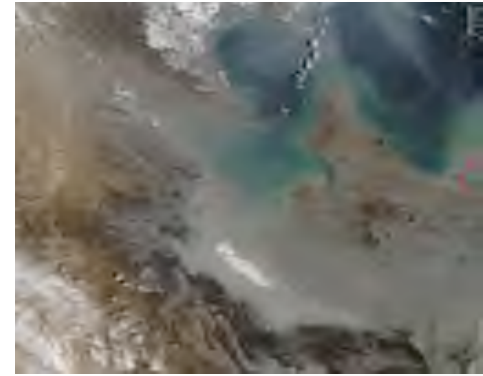
Smoke– True Positive



Hurricane – False Negative



Dust – False Positive



Smoke– False  
Positive



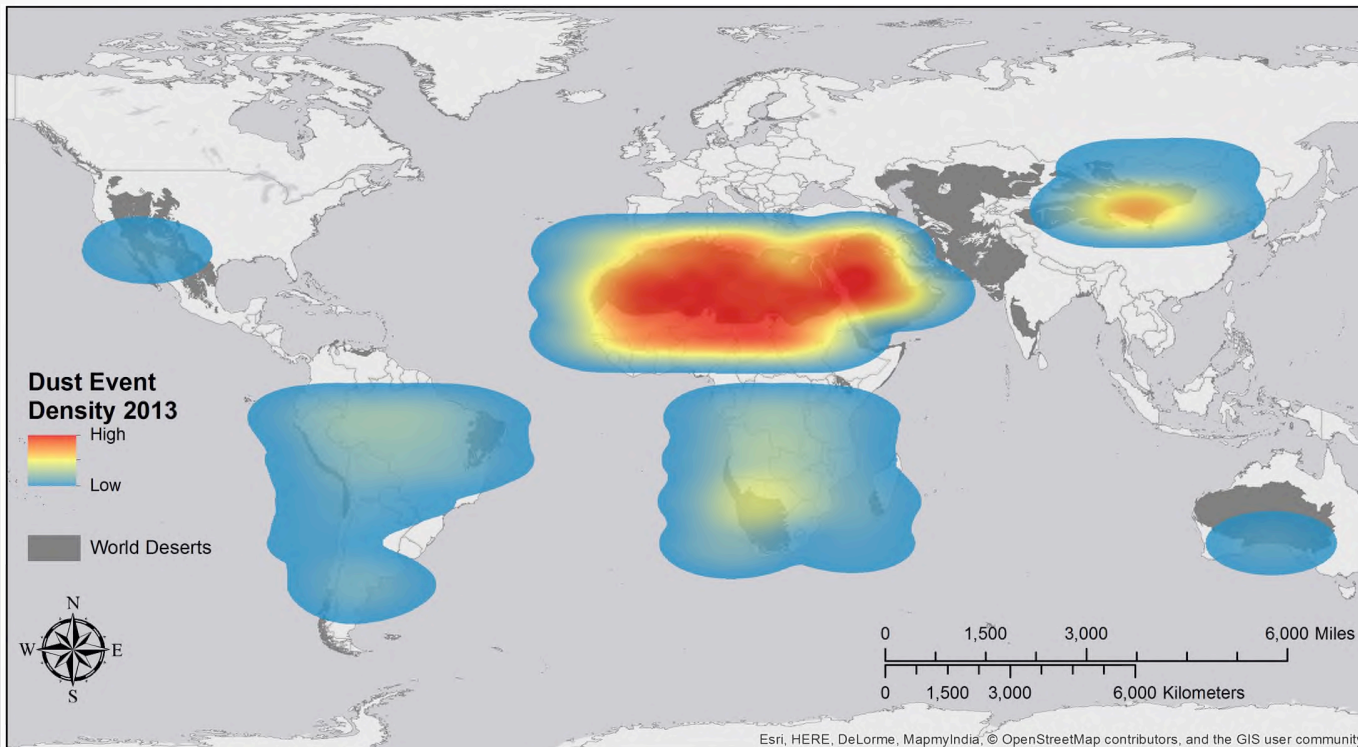
# Applications: Enabling new science

- Dust climatology – Collaboration with Sundar Christopher, UAH Atmospheric Science Professor

True\Predicted	Dust	Other	Total
Dust	1379	379	1758
Other	260	4932	5192
	1639	5311	6950

Validation  
Accuracy = **91%**

**Confusion Matrix**

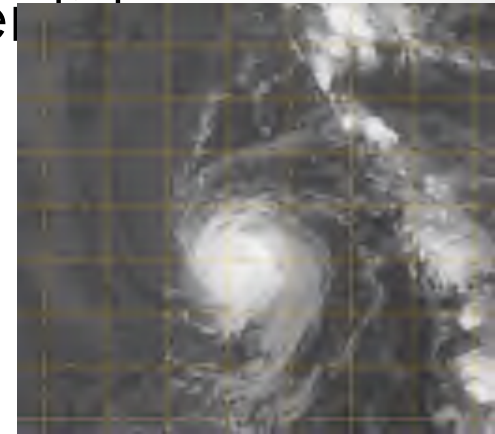


**Based on GIBS**

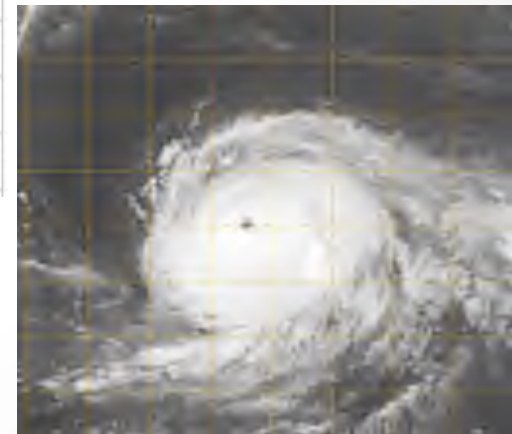
# Applications: Improving forecast operations

- Hurricane intensity estimation - Collaboration with Dan Cecil, NASA/MSFC Atmospheric Sciences

True\Predicted	td	ts	h1	h2	h3	h4	h5	no_cat	total
td	3168	335	0	1	0	0	0	6	3510
ts	489	4823	159	5	11	3	6	0	5496
h1	9	484	1158	92	20	6	1	0	1770
h2	3	76	214	513	145	4	0	5	960
h3	6	40	33	155	689	55	0	0	978
h4	1	18	17	12	142	810	32	0	1032
h5	2	2	0	0	27	59	216	0	306
no_cat	22	0	0	0	0	0	0	32	54
	3700	5778	1581	778	1034	937	255	43	14106



**Cat 2 Hurricane**

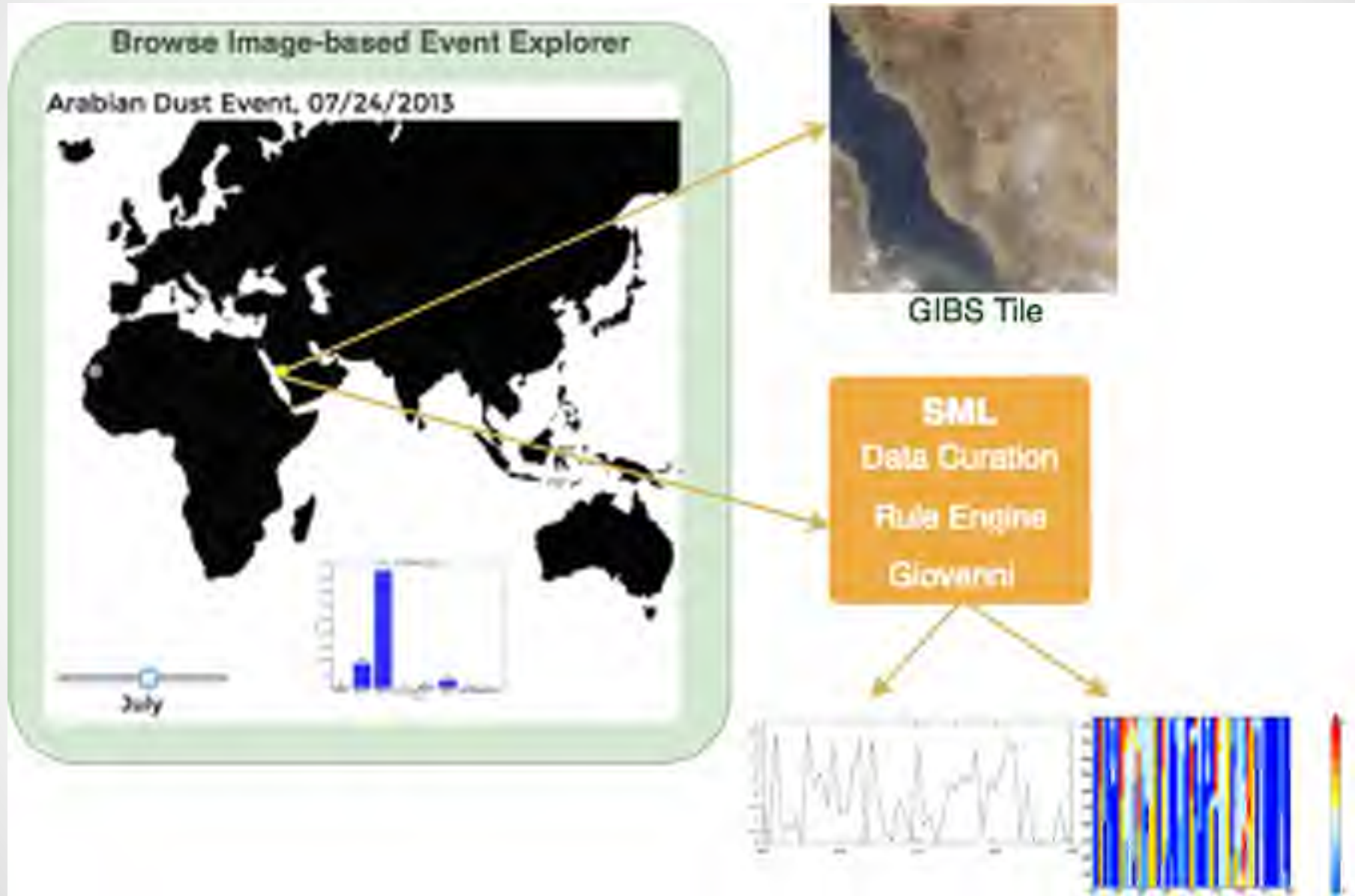


**Cat 4 Hurricane**

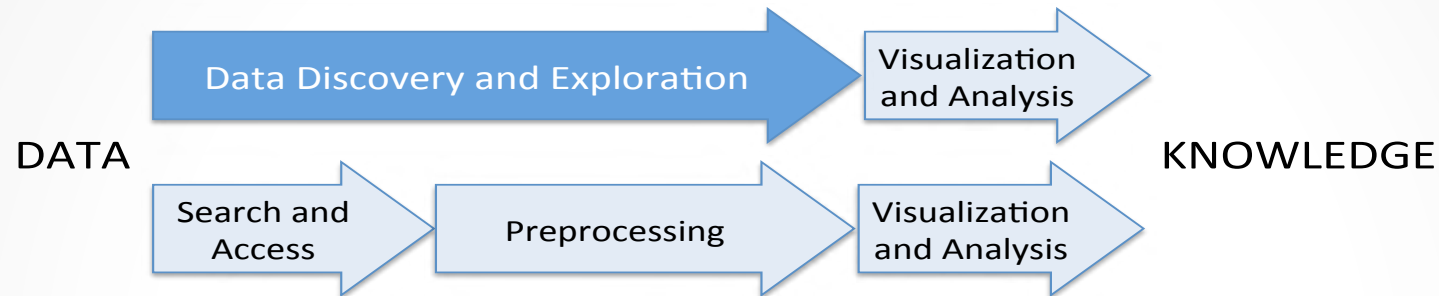
**Overall Accuracy : 81 %(Top 2 Probabilities 95.73%)**

**Data: NRL Images, HURDAT**

# Ongoing Work



# Summary



- Science data and information systems need to evolve to enable better data *search, access and usability!*
- Need operational services like – Data Curation Service, Rules Engine and Image Retrieval

# Questions

**Dr. Rahul Ramachandran**

Deputy Editor, Earth Science Informatics Journal

Manager, Global Hydrology Resource Center

NASA's Distributed Active Archive Center

Earth Science Office (ZP11)

NASA / Marshall Space Flight Center

Huntsville, Alabama 35812, USA

(w) 256.961.7620

(c) 256.226.6854

Orcid: [orcid.org/0000-0002-0647-1941](https://orcid.org/0000-0002-0647-1941)

[www.linkedin.com/in/ramachandran05/](https://www.linkedin.com/in/ramachandran05/)